

UNITED STATES PATENT APPLICATION

**SYSTEMS AND METHODS FOR ROUTING PACKETS
IN MULTIPROCESSOR COMPUTER SYSTEMS**

INVENTOR

Paul Krueger
of Minnetonka, Minnesota, USA

Schwegman, Lundberg, Woessner, & Kluth, P.A.
1600 TCF Tower
121 South Eighth Street
Minneapolis, Minnesota 55402
ATTORNEY DOCKET 1376.727US1

SYSTEMS AND METHODS FOR ROUTING PACKETS IN MULTIPROCESSOR COMPUTER SYSTEMS

5 Field of the Invention

The present invention is relates generally to the field of high-speed digital data processing systems, and more particularly to systems and methods of routing packets in multiprocessor computer systems.

Background Information

10 Multiprocessor computer systems consist of nodes interconnected by physical communication links in an n-dimensional topology. Messages between nodes are routed across the physical communication links in a variety of ways. In one approach, such as is described in U.S. Patent No. 5,970,232, issued October 19, 1999 to Passint et al., routers route messages between pairs of processing element nodes and a three
15 dimensional network. A symmetric eight port router is described that, in one embodiment, has a port connected to each node; the remaining ports are connected to physical communication links in the +X, -X, +Y, -Y, +Z and -Z directions.

A primary consideration in the design of interconnect networks and corresponding routing algorithms is avoiding deadlock. Deadlock occurs when cyclic
20 dependencies arise among a set of channel buffers, causing all involved buffers to fill up and block.

Approaches for avoiding deadlock are often dependent on a regular topology; failure of one or more communication links can reintroduce cyclic dependencies into what had been a deadlock-free routing scheme. What is needed is a system and method
25 of routing that avoids cyclic dependencies in networks with irregular topologies, and in regular topologies made irregular due to failures on one or more components.

Brief Description of the Drawings

Fig. 1 illustrates a multiprocessor computer system according to the present invention;

Fig. 2 illustrates a processing node which can be used in the multiprocessor computer system of Fig. 1;

Fig. 3 illustrates a method of building a routing table according to the present invention; and

Fig. 4 shows another example of a multiprocessor computer system according to the present invention.

10

Description of the Preferred Embodiments

In the following detailed description of the preferred embodiments, reference is made to the accompanying drawings which form a part hereof, and in which is shown by way of illustration specific embodiments in which the invention may be practiced. It is to be understood that other embodiments may be utilized and structural changes may be made without departing from the scope of the present invention.

A representative multiprocessor system 10 is shown in Fig. 1. In the example shown in Fig. 1, processing nodes 12 are connected to each other via communication links 14. In the embodiment shown, each node is connected directly to four other nodes. The present approach can be used on other topologies as well.

A representative processing node 12 is shown in Fig. 2. In the example shown in Fig. 2, a processing node includes a processor 20, a router 22 and memory 24. Router 22 includes a memory port 26, a processor port 28, four router ports 30 and a routing table 32. Router ports 30 connected to other nodes 12 over links 14.

In one embodiment, all routing across links 14 is controlled by software-loaded, distributed routing tables. This allows flexibility in the topology and network configuration. In one embodiment, routing table 32 in router 22 specifies the router port

30 to take based upon the destination node number and address. In one such embodiment, a separate routing table 32 is associated with each router port 30.

5 In one such embodiment, system 10 provides global addressing of local memories 24. High-order address bits specify the node 12 that is the destination of the message.

10 In one embodiment, routing tables 32 are consulted in parallel with arbitration, and actually specify the router ports 30 of the next router 22, so that the desired outputs can be known as soon as a packet enters the next router 22. In addition, each router 22 includes a buffer pool at each input for each virtual channel. In one embodiment, a separate virtual channel is provided for request and response packets. All packets are eligible to arbitrate for output ports (packets are not blocked behind packets ahead of them on the same virtual channel). Router 22 implements a wave-based arbitration mechanism that allows the highest-priority packet in arbitration to choose its output first, and then assigns priority for each port 30 in a round-robin manner.

15 In one embodiment, packet priority is selectable by software, such that certain classes of packets (e.g., I/O) can be given higher priority. The network also automatically ages packets in the network, increasing their priority the longer they stay in the network. This reduces the incidence of extremely long wait times that can otherwise result from multiple stages of contention within the network.

20 Each routing table 32 is created incrementally. Smaller routes are developed and then used to build larger routes. The algorithm applies a form of greedy algorithm to converge quickly to a functionally correct, but perhaps sub-optimal result. One example of building a routing table according to the present invention is described in Fig. 3.

25 In the embodiment shown in Fig. 3, the node initializes at 60 and each node 12 determines its "0" routes (i.e., the routes to itself). Control then moves to 62, where each node 12 queries its neighboring nodes 12 for their "0" hops. These are used by each node to determine all "1" hops from that node. A check is made at 66 if there is more than one path to any node. If so, all routes that are longer than the route of the

shortest route (or routes) are discarded. A check is made at 68 to determine if any of the new routes introduce a cycle. If so, the route is discarded.

Control then moves to 69, where the node selects among all the routes of length i that haven't been eliminated by the previous steps. Those preferred choices become the routes that are given to its neighbors in step 62 of the next iteration (route length of $X+1$).

In one embodiment, each node 12 reviews the list of routes at 69. If the list of routes has more than one route to any particular node, the node applies a routing algorithm to pick the route used in routing to that node. In one embodiment, a greedy algorithm is used to select the route based on a routing protocol such as dimension order routing. In another embodiment, each node looks at the number of times a particular link 14 is traversed by its routes, and attempts to spread the load across all the possible links 14.

Control then moves to 70. At 70, each node 12 determines if all other nodes can be reached by a route. If so, control moves to 72. If not, control moves to 70 and increments variable X .

Control then moves to 62, where each node 12 queries its neighboring nodes 12 for their " X " hops. These are used by each node 12 to determine all " $X + 1$ " hops from that node. A check is made at 66 if there is more than one path to any node. If so, all routes that are longer than the route of the shortest route (or routes) are discarded.

A check is made at 68 to determine if any of the new routes introduce a cycle. If so, the route is discarded. At 70, each node 12 determines if all other nodes can be reached by a route. If so, control moves to 74. If not, control moves to 72 and increments variable i .

At 74, each node builds its routing table.

Deadlock can also arise due to dependencies between request and response packets. Cycles are introduced between request and response packets if they use

dependent resources. One mechanism used to avoid dependency cycles is the use of alternative virtual channels.

In one embodiment, system 10 includes two virtual channels for each communication link 14. In one such embodiment, the hardware provides for requests
5 and responses to use different virtual channels.

In addition, the routing tables can be used to select from two virtual channels for each request or response. Each possible virtual channel selection becomes just another route choice for the algorithm. Other routing schemes make explicit choices about the use of virtual channels (for example the use of explicitly pre-selected "datelines"). In
10 this algorithm no special manipulation of virtual channels is required. Instead, the use of alternative virtual channels to avoid cycles becomes an emergent property of the algorithm. In a normally configured network the result may be identical to that of schemes that use explicit datelines. In an irregular topology, such as those that arise as the result of faults in hardware, they may be used differently. In all cases, the algorithm
15 assures routes that are free of dependency cycles.

Deadlock situations can be shown via a channel dependency graph, a directed graph whose nodes represent network channels and whose arcs represent dependencies between channels. An arc exists between channels x and y if and only if a packet will be routed directly from channel x to channel y as part of some selected route between
20 two nodes. It can be proven that a network is deadlock free if its channel dependency graph is acyclic.

To avoid deadlock due to cyclic dependency, in one embodiment a dateline is defined. Any packet that crosses that dateline changes from one virtual channel to another. That breaks the dependency cycle.

25 Such an approach works well for regular topologies but not as well for irregular topologies. In one embodiment, therefore, each node 12 combats cyclic dependency by maintaining a bit vector for each virtual channel. They are used as part of 68 to eliminate cycle routes.

A bit-vector is kept for each hop (call it "A") (a hop is defined as the use of a specific virtual channel over a specific physical link). Each bit in the bit vector represents some other hop in the system. A bit corresponding to hop "B" is set in the bit-vector for hop "A" if and only if the inclusion of a dependency from hop A to hop B
5 would result in a cycle in the dependency graph. The specific data structures and algorithm used to support this are discussed below. The result is that any dependency cycle that would be introduced as the result of adding any route can be checked by looking at a single bit.

Since the routing algorithm always builds longer routes from shorter routes, it is
10 sufficient to check only the single dependency that arises from the first two hops of any given route. All subsequent dependencies that arise because of routes of length i would have been checked when shorter sub-routes were selected in previous stages of the algorithm. When a route is selected, all relevant bit-vectors are updated for use in subsequent stages of the routing algorithm.

15 In one such embodiment, each hop is described in terms of the source and destination output port 30 and the virtual channel being used. Each instance of a hop must have a unique reference. That is, if it is necessary to refer to the same <port port vc> combination in multiple places, those reference must be to the same physical structure. This can be done in code by providing a function that either locates an
20 existing hop structure with specified ports and virtual channel or else creates such a structure if it doesn't already exist.

One way to implement such a function is to keep a global array of hops and search it as needed. It should be sized to be large enough to handle as many unique hops as are possible within the network being routed. For example, if there are four inter-
25 router links leaving each router 22, with two possible virtual channels defined for each port 30, the number of possible hops is $8 * \text{the number of routers}$ 30.

In one approach, one could use the index of the hop's entry in the array as a hop identifier and have a field in the hop itself that contains this value. One could also keep

a value in the system that points to the next unallocated entry in the array and update this as new hops are created.

At the beginning of the routing algorithm, the global hop-array is cleared. Each hop object includes an index and four bit-vectors (boolean arrays) that represent parts of the dependency graph and are used as part of the cycle checking. In the description
5 below, the four bit-vectors will be referred to as: cycle-to-hops, cycle-from-hops, pred-hops and succ-hops. These are sized identically to the global hop array with one bit for each possible hop. The dependency graph is consequently kept in a distributed fashion among all the hops.

10 When a dependency is added between, for example, hopA and hopB, the bit corresponding to hopB is set in the succ-hops array of hopA and the bit corresponding to hopA is set in the pred-hops array of hopB. Other updates are required as shown below. (The following example assumes the presence of functions to do an "inclusive or" of two arrays (bit-ior that takes two arrays as arguments and leaves the result in the
15 first array) and to set individual bits (sbit that takes an array and an index to set and a bit value)).

A pseudo code representation of the actions to be taken when a dependency is added (i.e. we have accepted a new route and add the dependency between the first two hops in the route) follows.

20 1) Update hops that would result in a cycle when going from hopB:

bit-ior (hopB->cycle-from-hops, hopA->cycle-from-hops)

bit-ior (hopB->cycle-from-hops, hopA->pred-hops)

2) Update hops that would result in a cycle when going to hopA:

bit-ior (hopA->cycle-to-hops, hopB->cycle-to-hops)

25 bit-ior (hopA->cycle-to-hops, hopB->succ-hops)

3) Update pred-hops for hopB:

sbit (hopB->pred-hops, hopA->index, 1)

4) Update succ-hops for hopA

```
sbit (hopA->succ-hops, hopB->index, 1)
```

5) Update hops that would result in a cycle when going from some other hop to hopA.

(Assume hops is the global system array of hops. You need to iterate over a set of hop

5 indices. You can do this using a function that converts a bit array to a list of corresponding numbers (i.e. if bit X is set then the value X appears in the output list).)

```
it_list = bit-vector-to-numbers (hopA->cycle-to-hops)
```

```
for (i=it_list.first, i!=null, i=it_list.next)
```

```
sbit (hops[i]->cycle-from-hops, hopA->index, 1)
```

10 bit-ior (hops[i]->cycle-from-hops, hopA->pred-hops)

```
bit-ior (hops[i]->cycle-from-hops, hopA->cycle-from-hops)
```

6) Update hops that would result in a cycle when going from some other hop to hop B:

```
it_list = bit-vector-to-numbers (hopB->cycle-from-hops)
```

```
for (i=it_list.first, i!=null, i=it_list.next)
```

15 sbit (hops[i]->cycle-to-hops, hopB->index 1)

```
bit-ior (hops[i]->cycle-to-hops, hopB->succ-hops)
```

```
bit-ior (hops[i]->cycle-to-hops, hopB->cycle-to-hops)
```

While routing, to check for a cycle if the dependency hopA=>hopB were to be added to the graph, use the following check:

20 if (hopA->cycle-from-hops [hopB->index] == 1) ...

If the condition is met, then a cycle would exist if that dependency was added and the route will be rejected. By making the incremental algorithm that constructs the routing table aware of virtual channels and by allowing each node 12 to choose between virtual channels, the algorithm in effect defines its own date line equivalents wherever it needs them. They arise naturally as an emergent property of the algorithm rather than

25 having been explicitly defined. When there is only one available route, it is taken.

When, however, there are many possible routes, node 12 applies a greedy algorithm to choose between the routes.

The general notion of a greedy algorithm is that, given a series of decisions to make, they should be made as quickly and optimally as possible. This may result in
5 choices that are not globally optimal. However, in most cases using reasonable decision heuristics results in choices that are close to or exactly globally optimal. The idea behind greedy is that each decision is the best one that can immediately be made using currently available information. Choices are never re-considered or changed. Greedy algorithms, therefore, tend to be very quick to make decisions. Non-optimality can be
10 mitigated to some extent by providing very good guidelines for making the needed decisions.

For example, one of the standard ways of routing is dimension order route. This choice criterion can be used within this invention to select from among alternative routes that one which most closely follows a specified dimension order. If applying this
15 to a regular, fault-free topology, the described algorithm of the present invention will always choose to move in a designated dimension first, and then in a second dimension. The set of selected routes will consequently be globally optimal.

When, however, there are faults (e.g. a missing link between two nodes) or the dimension ordered routing is not possible for some other reason, the routing algorithm
20 of the present invention is able to adapt to those faults, deviate from strict dimension ordering, and find an alternative route while guaranteeing that there are no cyclical dependencies.

Another example of a multiprocessor system 10 is shown in Fig. 4. In the example shown in Fig. 4, a router 30 in each processing node 12 is connected to a router
25 16 and to a router 30 in another node 12. Each router 16 includes eight router ports 18. Four ports 18 are connected to processing nodes 12; the remaining ports 18 are connected to other routers. In one embodiment, processing node 12 is configured as shown in Fig. 2.

As in the example discussed above, in one embodiment of system 10 in Fig. 4, all routing across links 14 is controlled by software-loaded, distributed routing tables. A router table is stored in memory of each router 16 and each router 30. This allows flexibility in the topology and network configuration. In one embodiment, each routing table in router 16 specifies the router port 18 to take based upon the destination node number and address.

The method of incrementally building a routing table applies to the configuration shown in Fig. 4 as well. Once again, routing tables are built incrementally by querying neighbors as in Fig. 3. Now, however, router 16 becomes a neighbor in the process. The method of identifying cyclic dependencies applies as well.

Now, however, hops will be used for router to node connections as well. That will increase the size of the array that is needed.

Another example of a multiprocessor system 10 is shown in Fig. 4. In the example shown in Fig. 4, a router 30 in each processing node 12 is connected to a router 16 and to a router 30 in another node 12. Each router 16 includes eight router ports 18. Four ports 18 are connected to processing nodes 12; the remaining ports 18 are connected to other routers. In one embodiment, processing node 12 is configured as shown in Fig. 2.

As in the example discussed above, in one embodiment of system 10 in Fig. 4, all routing across links 14 is controlled by software-loaded, distributed routing tables. A router table is stored in memory of each router 16 and each router 30. This allows flexibility in the topology and network configuration. In one embodiment, each routing table in router 16 specifies the router port 18 to take based upon the destination node number and address.

The method of incrementally building a routing table applies to the configuration shown in Fig. 4 as well. Once again, routing tables are built incrementally by querying neighbors as in Fig. 3. Now, however, router 16 becomes a neighbor in the process. The method of identifying cyclic dependencies applies as well. Now,

however, hops will be used for router to node connections as well. That will increase the size of the array you need.

Definitions

5 In the above discussion, the term “computer” is defined to include any digital or analog data processing unit. Examples include any personal computer, workstation, set top box, mainframe, server, supercomputer, laptop or personal digital assistant capable of embodying the inventions described herein.

10 Examples of articles comprising computer readable media are floppy disks, hard drives, CD-ROM or DVD media or any other read-write or read-only memory device.

 Portions of the above description have been presented in terms of algorithms and symbolic representations of operations on data bits within a computer memory. These algorithmic descriptions and representations are the ways used by those skilled in the data processing arts to most effectively convey the substance of their work to others skilled in the art. An algorithm is here, and generally, conceived to be a self-consistent sequence of steps leading to a desired result. The steps are those requiring physical manipulations of physical quantities. Usually, though not necessarily, these quantities take the form of electrical or magnetic signals capable of being stored, transferred, combined, compared, and otherwise manipulated. It has proven convenient at times, 15 principally for reasons of common usage, to refer to these signals as bits, values, elements, symbols, characters, terms, numbers, or the like. It should be borne in mind, however, that all of these and similar terms are to be associated with the appropriate physical quantities and are merely convenient labels applied to these quantities. Unless specifically stated otherwise as apparent from the following discussions, terms such as 20 “processing” or “computing” or “calculating” or “determining” or “displaying” or the like, refer to the action and processes of a computer system, or similar computing device, that manipulates and transforms data represented as physical (e.g., electronic) quantities within the computer system=s registers and memories into other data

similarly represented as physical quantities within the computer system memories or registers or other such information storage, transmission or display devices.

Although specific embodiments have been illustrated and described herein, it will be appreciated by those of ordinary skill in the art that any arrangement which is
5 calculated to achieve the same purpose may be substituted for the specific embodiment shown. This application is intended to cover any adaptations or variations of the present invention. Therefore, it is intended that this invention be limited only by the claims and the equivalents thereof.